

# Convex Optimization

## CMU-10725

Independent Component Analysis  
(and matrix differentials)

Barnabás Póczos & Ryan Tibshirani



**MACHINE LEARNING** DEPARTMENT

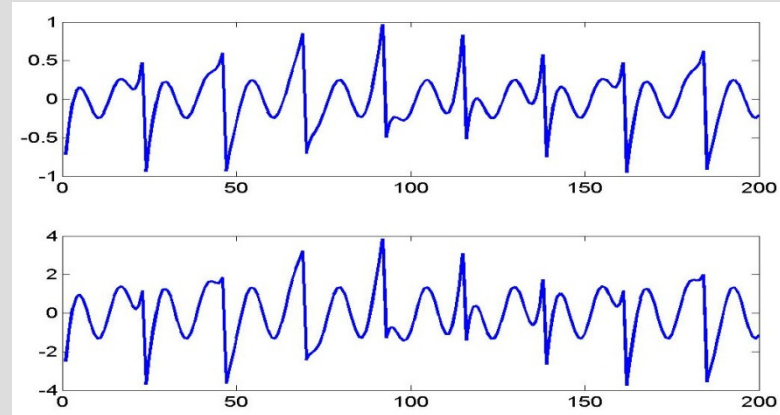


# Independent Component Analysis

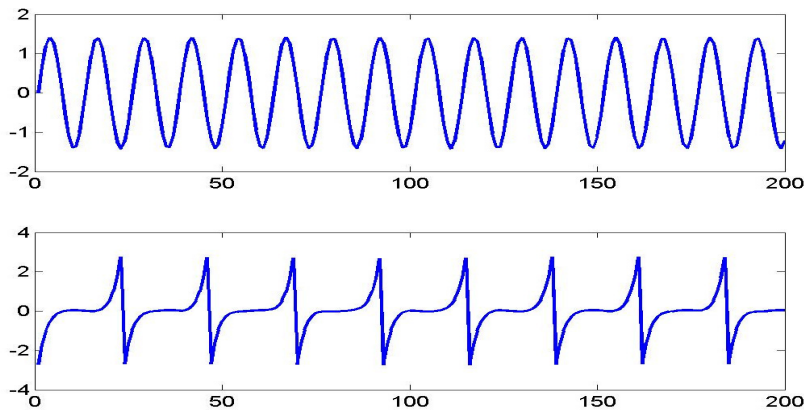
# Independent Component Analysis

$$\begin{aligned}x_1(t) &= a_{11}s_1(t) + a_{12}s_2(t) \\x_2(t) &= a_{21}s_1(t) + a_{22}s_2(t)\end{aligned}$$

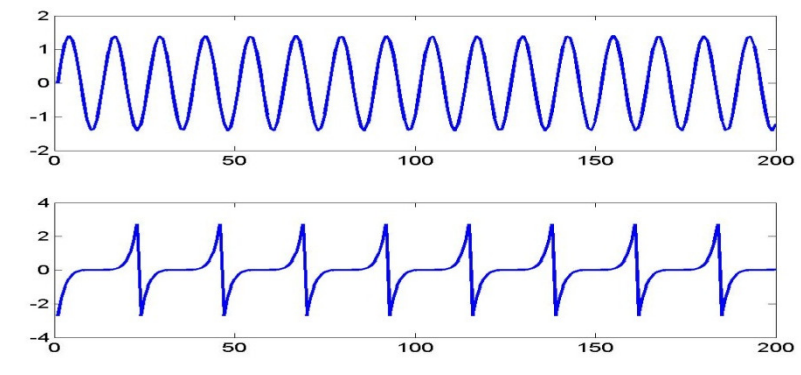
Model



Observations (Mixtures)



ICA estimated signals



original signals

# Independent Component Analysis

**Model**

$$x_1(t) = a_{11}s_1(t) + a_{12}s_2(t)$$

$$x_2(t) = a_{21}s_1(t) + a_{22}s_2(t)$$

**We observe**

$$\begin{pmatrix} x_1(1) \\ x_2(1) \end{pmatrix}, \begin{pmatrix} x_1(2) \\ x_2(2) \end{pmatrix}, \dots, \begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix}$$

**We want**

$$\begin{pmatrix} s_1(1) \\ s_2(1) \end{pmatrix}, \begin{pmatrix} s_1(2) \\ s_2(2) \end{pmatrix}, \dots, \begin{pmatrix} s_1(t) \\ s_2(t) \end{pmatrix}$$

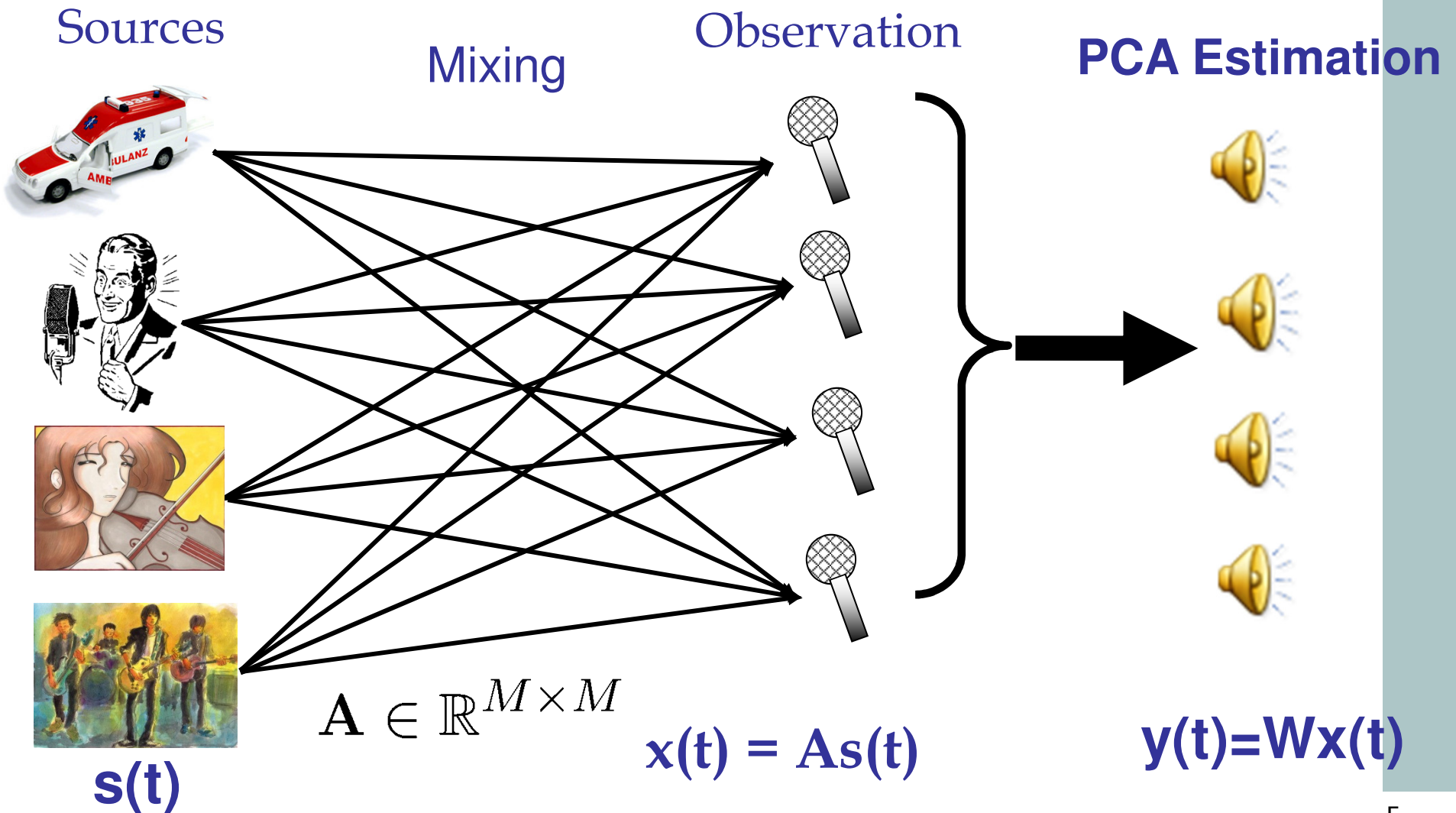
But we don't know  $\{a_{ij}\}$ , nor  $\{s_i(t)\}$

**Goal:**

Estimate  $\{s_i(t)\}$ , (and also  $\{a_{ij}\}$ )

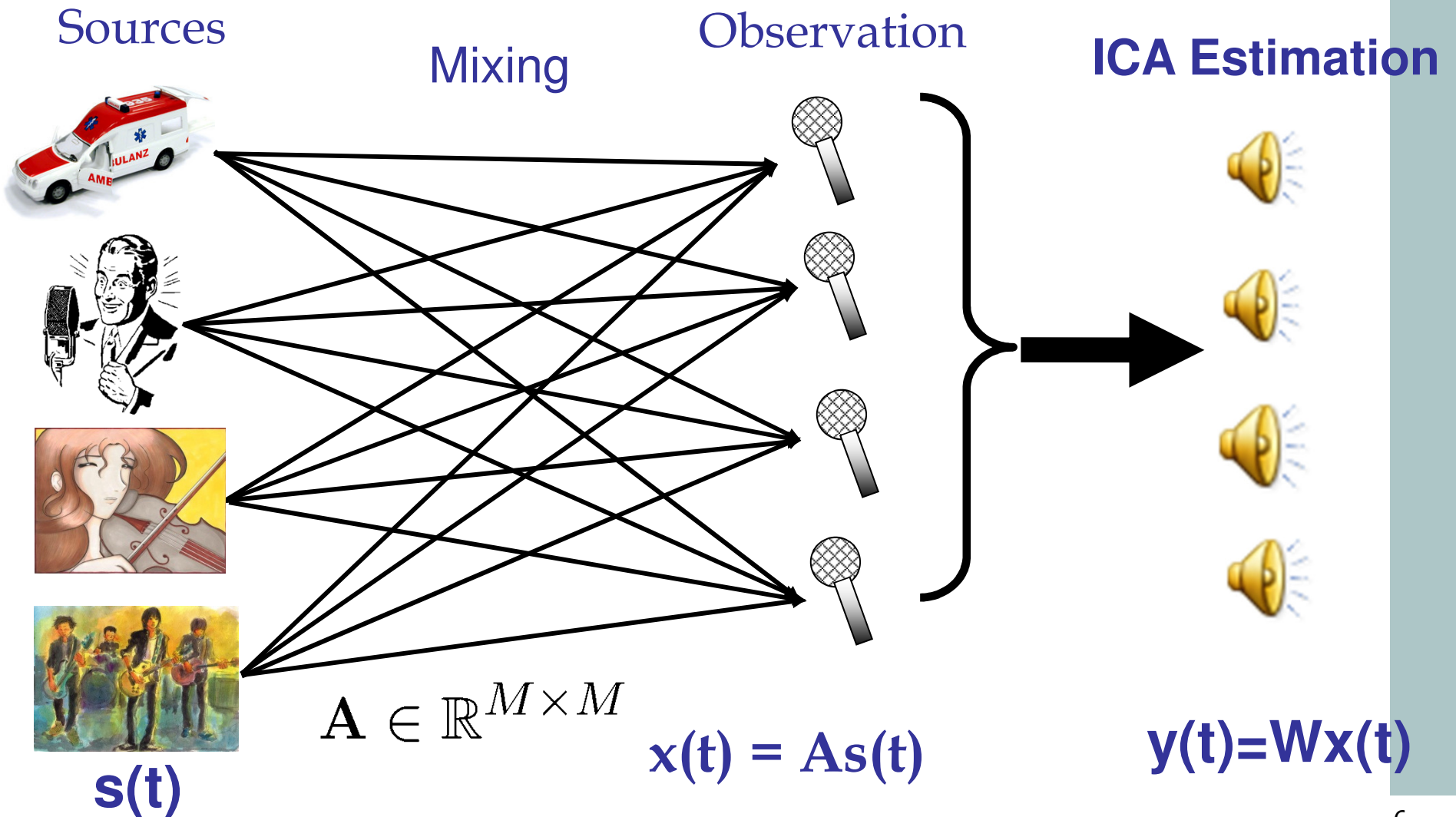
# The Cocktail Party Problem

## SOLVING WITH PCA



# The Cocktail Party Problem

## SOLVING WITH ICA



# ICA Theory

# Statistical (in)dependence

## Definition (Independence)

$Y_1, Y_2$  are independent  $\Leftrightarrow p(y_1, y_2) = p(y_1)p(y_2)$

## Definition (Shannon entropy)

$$H(\mathbf{Y}) \doteq H(Y_1, \dots, Y_m) \doteq - \int p(y_1, \dots, y_m) \log p(y_1, \dots, y_m) d\mathbf{y}.$$

## Definition (KL divergence)

$$0 \leq KL(f||g) = \int f(x) \log \frac{f(x)}{g(x)} dx$$

## Definition (Mutual Information)

$$0 \leq I(Y_1, \dots, Y_M) \doteq \int p(y_1, \dots, y_M) \log \frac{p(y_1, \dots, y_M)}{p(y_1) \dots p(y_M)} d\mathbf{y} \quad 8$$



# Solving the ICA problem with i.i.d. sources

**ICA problem:**  $\mathbf{x} = \mathbf{A}\mathbf{s}$ ,  $\mathbf{s} = [s_1; \dots; s_M]$  are jointly independent.

## Ambiguity:

$\mathbf{s} = [s_1; \dots; s_M]$  sources can be recovered only up to  
**sign, scale and permutation.**

## Proof:

- $\mathbf{P}$  = arbitrary permutation matrix,
- $\mathbf{\Lambda}$  = arbitrary diagonal scaling matrix.

$$\Rightarrow \mathbf{x} = [\mathbf{A}\mathbf{P}^{-1}\mathbf{\Lambda}^{-1}][\mathbf{\Lambda}\mathbf{P}\mathbf{s}]$$

# Solving the ICA problem

## Lemma:

We can assume that  $E[s] = 0$ .

## Proof:

Removing the mean does not change the mixing matrix.

$$\mathbf{x} - E[\mathbf{x}] = \mathbf{A}(\mathbf{s} - E[\mathbf{s}]).$$

In what follows we assume that  $E[\mathbf{s}\mathbf{s}^T] = \mathbf{I}_M$ ,  $E[\mathbf{s}] = 0$ .

# Whitening

- Let  $\Sigma \doteq \text{cov}(\mathbf{x}) = E[\mathbf{x}\mathbf{x}^T] = \mathbf{A}E[\mathbf{s}\mathbf{s}^T]\mathbf{A}^T = \mathbf{A}\mathbf{A}^T$ .  
(We assumed centered data)
- Do **SVD**:  $\Sigma \in \mathbb{R}^{N \times N}$ ,  $\text{rank}(\Sigma) = M$ ,  
 $\Rightarrow \Sigma = \mathbf{U}\mathbf{D}\mathbf{U}^T$ ,  
where  $\mathbf{U} \in \mathbb{R}^{N \times M}$ ,  $\mathbf{U}^T\mathbf{U} = \mathbf{I}_M$ , **Singular vectors**  
 $\mathbf{D} \in \mathbb{R}^{M \times M}$ , diagonal with rank  $M$ . **Singular values**

# Whitening

- Let  $Q \doteq D^{-1/2}U^T \in \mathbb{R}^{M \times N}$  *whitening matrix*
- Let  $A^* \doteq QA$
- $x^* \doteq Qx = QA s = A^* s$  is our new (*whitened*) ICA task.

**We have,**

$$E[x^* x^{*T}] = E[Qxx^T Q^T] = Q\Sigma Q^T = (D^{-1/2}U^T)UDU^T(UD^{-1/2}) = I_M$$

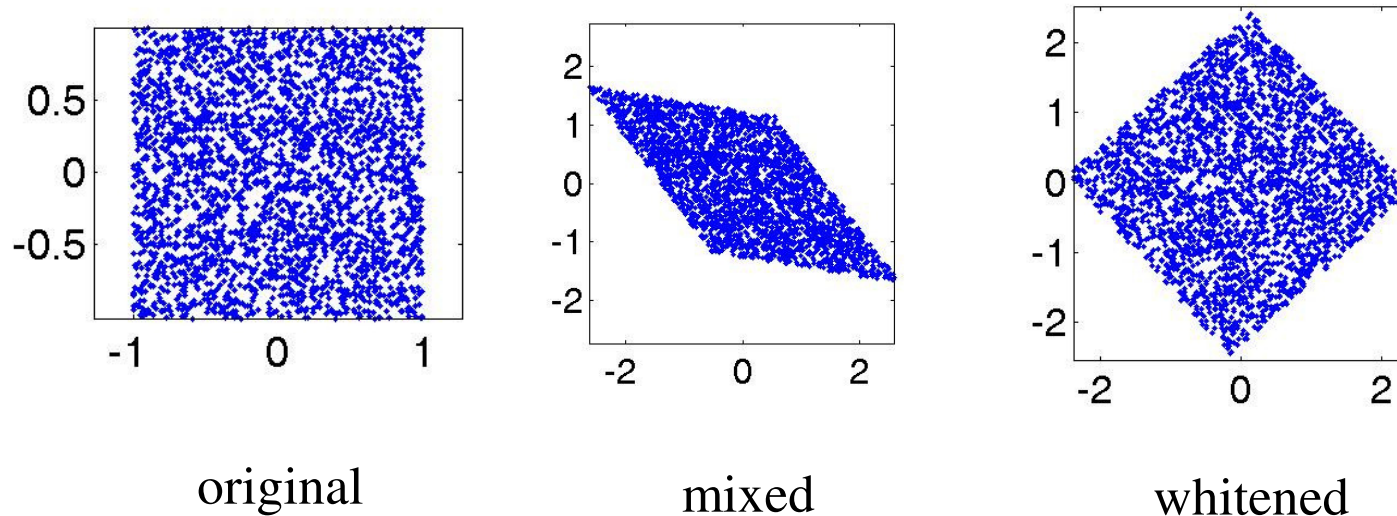
$$\Rightarrow E[x^* x^{*T}] = I_M, \text{ and } A^* A^{*T} = I_M.$$

# Whitening solves half of the ICA problem

## Note:

The number of free parameters of an  $N$  by  $N$  orthogonal matrix is  $(N-1)(N-2)/2$ .

⇒ whitening solves **half** of the ICA problem



After whitening it is enough to consider  
**orthogonal matrices** for separation.

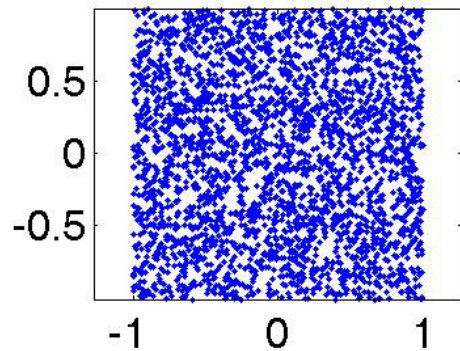
# Solving ICA

**ICA task:** Given  $\mathbf{x}$ ,

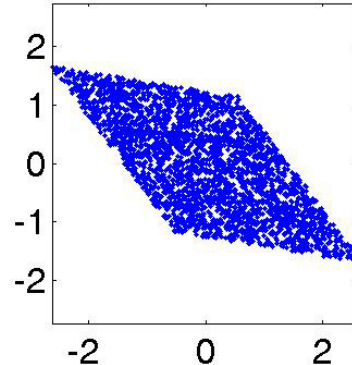
- ❑ find  $\mathbf{y}$  (the estimation of  $\mathbf{s}$ ),
- ❑ find  $\mathbf{W}$  (the estimation of  $\mathbf{A}^{-1}$ )

**ICA solution:**  $\mathbf{y} = \mathbf{W}\mathbf{x}$

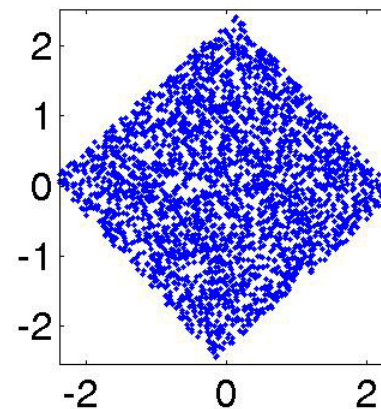
- ❑ Remove mean,  $E[\mathbf{x}] = 0$
- ❑ Whitening,  $E[\mathbf{x}\mathbf{x}^T] = \mathbf{I}$
- ❑ Find an orthogonal  $\mathbf{W}$  optimizing an objective function
  - Sequence of 2-d Jacobi (Givens) rotations



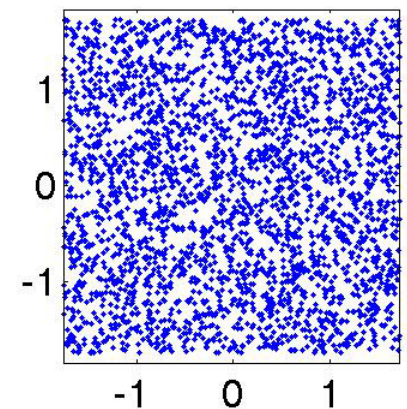
original



mixed



whitened



rotated

(demixed)

# Optimization Using Jacobi Rotation Matrices

$$G(p, q, \theta) \doteq \begin{pmatrix} 1 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \dots & \cos(\theta) & \dots & -\sin(\theta) & \dots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \dots & \sin(\theta) & \dots & \cos(\theta) & \dots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 1 \end{pmatrix} \begin{matrix} \leftarrow \mathbf{p} \\ \\ \leftarrow \mathbf{q} \end{matrix} \in \mathbf{R}^{M \times M}$$

$\uparrow$   
 $\mathbf{p}$

$\uparrow$   
 $\mathbf{q}$

*Observation* :  $\mathbf{x} = \mathbf{A}\mathbf{s}$

*Estimation* :  $\mathbf{y} = \mathbf{W}\mathbf{x}$

$$\mathbf{W} = \arg \min_{\tilde{\mathbf{W}} \in \mathcal{W}} J(\tilde{\mathbf{W}}\mathbf{x}),$$

$$\text{where } \mathcal{W} = \{\mathbf{W} | \mathbf{W} = \prod_i G(p_i, q_i, \theta_i)\}$$

# ICA Cost Functions

Let  $\mathbf{y} \doteq \mathbf{W}\mathbf{x}$ ,  $\mathbf{y} = [y_1; \dots; y_M]$ , and let us measure the dependence using Shannon's mutual information:

$$J_{ICA_1}(\mathbf{W}) \doteq I(y_1, \dots, y_M) \doteq \int p(y_1, \dots, y_M) \log \frac{p(y_1, \dots, y_M)}{p(y_1) \dots p(y_M)} d\mathbf{y},$$

Let  $H(\mathbf{y}) \doteq H(y_1, \dots, y_m) \doteq - \int p(y_1, \dots, y_m) \log p(y_1, \dots, y_m) d\mathbf{y}$ .

## Lemma

$$H(\mathbf{W}\mathbf{x}) = H(\mathbf{x}) + \log |\det \mathbf{W}| \quad \text{Proof: Homework}$$

Therefore,

$$\begin{aligned} I(y_1, \dots, y_M) &= \int p(y_1, \dots, y_M) \log \frac{p(y_1, \dots, y_M)}{p(y_1) \dots p(y_M)} \\ &= -H(y_1, \dots, y_M) + H(y_1) + \dots + H(y_M) \\ &= -H(x_1, \dots, x_M) - \log |\det \mathbf{W}| + H(y_1) + \dots + H(y_M). \end{aligned}$$



# ICA Cost Functions

$$\begin{aligned} I(y_1, \dots, y_M) &= \int p(y_1, \dots, y_M) \log \frac{p(y_1, \dots, y_M)}{p(y_1) \dots p(y_M)} \\ &= -H(y_1, \dots, y_M) + H(y_1) + \dots + H(y_M) \\ &= -H(x_1, \dots, x_M) - \log |\det \mathbf{W}| + H(y_1) + \dots + H(y_M). \end{aligned}$$

$H(x_1, \dots, x_M)$  is constant,  $\log |\det \mathbf{W}| = 0$ .

**Therefore,**

$$J_{ICA_2}(\mathbf{W}) \doteq H(y_1) + \dots + H(y_M)$$

The covariance is fixed: I. Which distribution has the largest entropy?

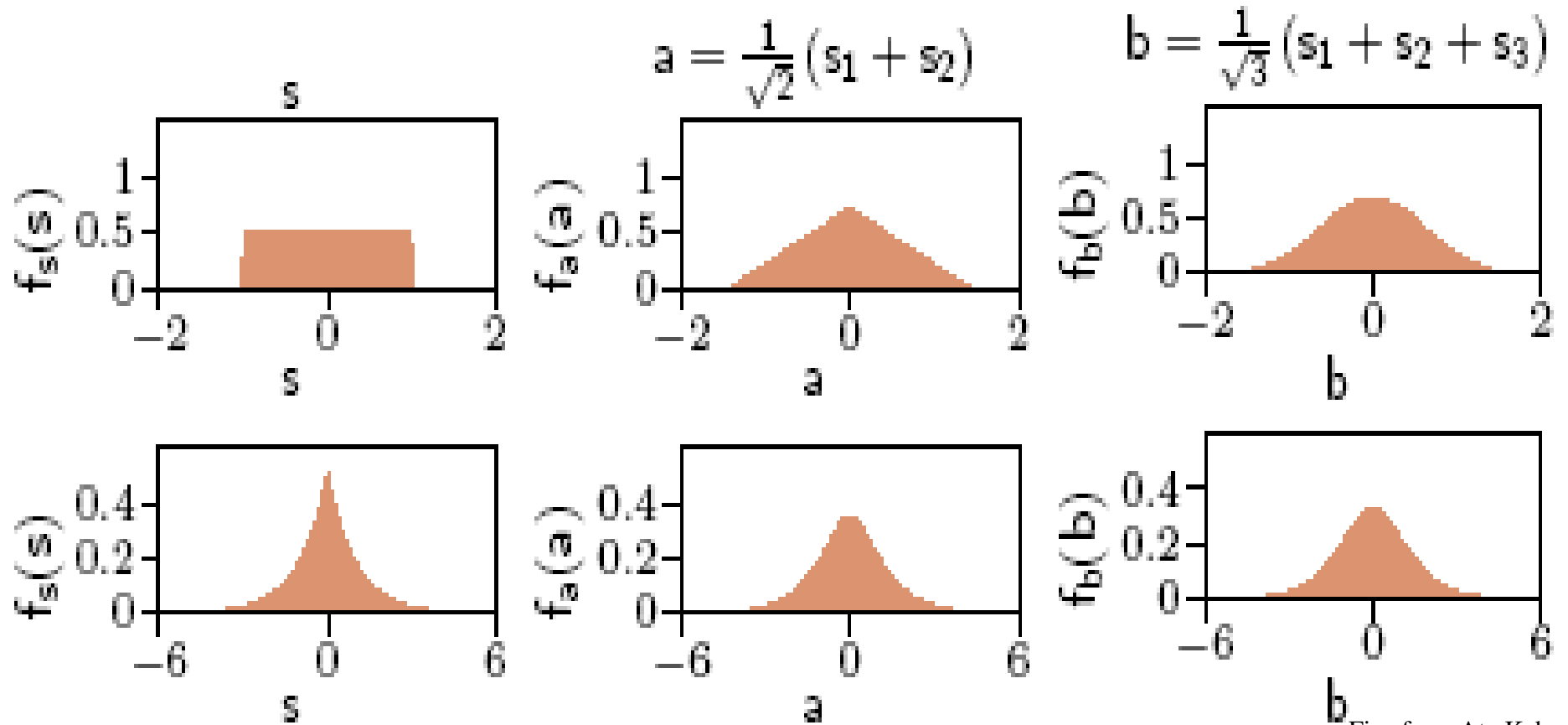
$\Rightarrow$  go away from normal distribution

# Central Limit Theorem

The sum of independent variables converges to the normal distribution

⇒ For separation go far away from the normal distribution

⇒ **Negentropy, |kurtosis| maximization**



# ICA Algorithms

# Maximum Likelihood ICA Algorithm

David J.C. MacKay (97)

- simplest approach
- requires knowing densities of hidden sources  $\{f_i\}$

rows of  $\mathbf{W}$

$\mathbf{x}(t) = \mathbf{A}\mathbf{s}(t)$ ,  $\mathbf{s}(t) = \mathbf{W}\mathbf{x}(t)$ , where  $\mathbf{A}^{-1} = \mathbf{W} = [\mathbf{w}_1; \dots; \mathbf{w}_M] \in \mathbb{R}^{M \times M}$

$$\begin{aligned}
 L &= \sum_{t=1}^T \log P_{\mathbf{x}}(\mathbf{x}(t)) = \sum_{t=1}^T \log \underbrace{P_{\mathbf{x}}(\mathbf{A}\mathbf{s}(t))}_{\mathbf{A}^{-1}P_{\mathbf{s}}(\mathbf{A}^{-1}\mathbf{x}(t))} \Rightarrow \max_{\mathbf{A}} \\
 P_{\mathbf{A}\mathbf{s}}(\mathbf{u}) &= \mathbf{A}^{-1}P_{\mathbf{s}}(\mathbf{A}^{-1}\mathbf{u}) \\
 \Rightarrow L &= \sum_{t=1}^T \log \mathbf{A}^{-1}P_{\mathbf{s}}(\mathbf{s}(t)) = T \log |\mathbf{W}| + \sum_{t=1}^T \log \underbrace{P_{\mathbf{s}}(\mathbf{s}(t))}_{\prod_{i=1}^M P_{s_i}(s_i(t))} \\
 &= T \log |\mathbf{W}| + \sum_{t=1}^T \sum_{i=1}^M \log \underbrace{P_{s_i}(\mathbf{w}_i^T \mathbf{x}(t))}_{f_i(\mathbf{w}_i^T \mathbf{x}(t))} \\
 &\Rightarrow \max_{\mathbf{W}}
 \end{aligned}$$

# Maximum Likelihood ICA Algorithm

$$L = T \log |W| + \sum_{t=1}^T \sum_{k=1}^M \log f_k(W_k x(t))$$

$$\Rightarrow \max_W L \Rightarrow \frac{\partial L}{\partial W_{ij}} = ?$$

$$\frac{\partial L}{\partial W_{ij}} = T(W^T)^{-1}_{ij} + \sum_{t=1}^T \underbrace{\frac{\partial}{\partial W_{ij}} \sum_{k=1}^M \log f_k(W_k x(t))}_{\frac{f'_i(W_i x(t))}{f_i(W_i x(t))} x_j(t)}$$

$$= T(W^T)^{-1}_{ij} + \sum_{t=1}^T \frac{f'_i(W_i x(t))}{f_i(W_i x(t))} x_j(t)$$

$$\Rightarrow \Delta W \propto [W^T]^{-1} + \frac{1}{T} \sum_{t=1}^T g(Wx(t)) x^T(t), \text{ where } g_i = f'_i/f_i$$

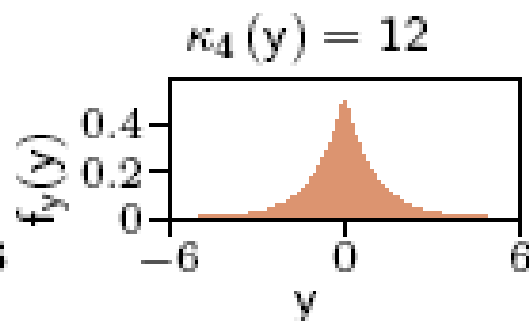
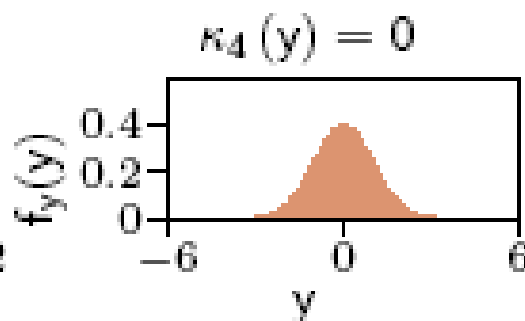
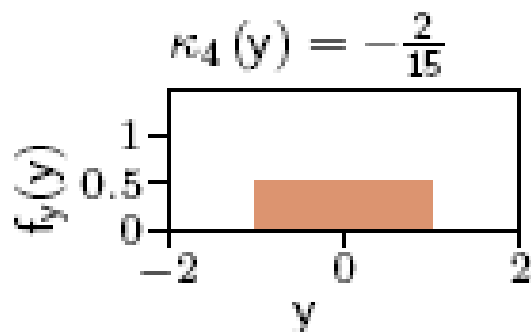
# ICA algorithm based on Kurtosis maximization

Kurtosis = 4<sup>th</sup> order cumulant

Measures

- the distance from normality
- the degree of peakedness

$$\bullet \kappa_4(y) = E\{y^4\} - \underbrace{3(E\{y^2\})^2}_{= 3 \text{ if } E\{y\} = 0 \text{ and whitened}}$$



# The Fast ICA algorithm (Hyvarinen)

- Given whitened data  $\mathbf{z}$
- Estimate the 1<sup>st</sup> ICA component:

Probably the most famous ICA algorithm

$$\star y = \mathbf{w}^T \mathbf{z}, \quad \|\mathbf{w}\| = 1, \quad \Leftarrow \mathbf{w}^T = 1^{st} \text{ row of } \mathbf{W}$$

$$\star \text{ maximize kurtosis } f(\mathbf{w}) \doteq \kappa_4(y) \doteq \mathbb{E}[y^4] - 3 \\ \text{with constraint } h(\mathbf{w}) = \|\mathbf{w}\|^2 - 1 = 0$$

$$\star \text{ At optimum } f'(\mathbf{w}) + \lambda h'(\mathbf{w}) = 0^T \quad (\lambda \text{ Lagrange multiplier})$$

$$\Rightarrow 4\mathbb{E}[(\mathbf{w}^T \mathbf{z})^3 \mathbf{z}] + 2\lambda \mathbf{w} = 0$$

Solve this equation by Newton–Raphson’s method.

# The Fast ICA algorithm (Hyvarinen)

**Solve:**  $F(\mathbf{w}) = 4\mathbb{E}[(\mathbf{w}^T \mathbf{z})^3 \mathbf{z}] + 2\lambda \mathbf{w} = 0$

**Note:**

$$y = \mathbf{w}^T \mathbf{z}, \quad \|\mathbf{w}\| = 1, \quad \mathbf{z} \text{ white} \Rightarrow \mathbb{E}[(\mathbf{w}^T \mathbf{z})^2] = 1$$

$$\mathbb{E}[\mathbf{z} \mathbf{z}^T] = \mathbf{I} \quad \xrightarrow{\quad} \quad \mathbb{E}[\mathbf{w}^T \mathbf{z} \mathbf{z}^T \mathbf{w}] = \mathbf{w}^T \mathbf{I} \mathbf{w} = 1$$

$\|\mathbf{w}\|=1$

**The derivative of  $F$ :**

$$\begin{aligned} F'(\mathbf{w}) &= 12\mathbb{E}[(\mathbf{w}^T \mathbf{z})^2 \mathbf{z} \mathbf{z}^T] + 2\lambda \mathbf{I} \\ &\sim 12\mathbb{E}[(\mathbf{w}^T \mathbf{z})^2] \mathbb{E}[\mathbf{z} \mathbf{z}^T] + 2\lambda \mathbf{I} \\ &= 12\mathbb{E}[(\mathbf{w}^T \mathbf{z})^2] \mathbf{I} + 2\lambda \mathbf{I} \\ &= 12\mathbf{I} + 2\lambda \mathbf{I} \end{aligned}$$



# The Fast ICA algorithm (Hyvarinen)

The Jacobian matrix becomes diagonal, and can easily be inverted.

$$\mathbf{w}(k+1) = \mathbf{w}(k) - [F'(\mathbf{w}(k))]^{-1} F(\mathbf{w}(k))$$

$$\mathbf{w}(k+1) = \mathbf{w}(k) - \frac{4\mathbb{E}[(\mathbf{w}(k)^T \mathbf{z})^3 \mathbf{z}] + 2\lambda \mathbf{w}(k)}{12 + 2\lambda}$$

$$(12 + 2\lambda)\mathbf{w}(k+1) = (12 + 2\lambda)\mathbf{w}(k) - 4\mathbb{E}[(\mathbf{w}(k)^T \mathbf{z})^3 \mathbf{z}] - 2\lambda \mathbf{w}(k)$$

$$-\frac{12+2\lambda}{4}\mathbf{w}(k+1) = -3\mathbf{w}(k) + \mathbb{E}[(\mathbf{w}(k)^T \mathbf{z})^3 \mathbf{z}]$$

Therefore,

Let  $\mathbf{w}_1$  be the fix pont of:

$$\tilde{\mathbf{w}}(k+1) = \mathbb{E}[(\mathbf{w}(k)^T \mathbf{z})^3 \mathbf{z}] - 3\mathbf{w}(k)$$

$$\mathbf{w}(k+1) = \frac{\tilde{\mathbf{w}}(k+1)}{\|\tilde{\mathbf{w}}(k+1)\|}$$

- Estimate the  $2^{nd}$  ICA component similarly  
using the  $\mathbf{w} \perp \mathbf{w}_1$  additional constraint... and so on ...

# Other Nonlinearities

$$\max_{\substack{\omega \\ \text{s.t. } \omega^T \omega = 1}} \mathbb{E} G(\omega^T z) \quad G(y) = y^4 - 3 \leftarrow \text{IN THE PREVIOUS}$$

$$\Rightarrow F(\omega) = \mathbb{E} z z^T g(\omega^T z) - \lambda \omega = 0 \in \mathbb{R}^n \quad \left\{ \begin{array}{l} g(y) = G'(y) \\ \quad \quad \quad \downarrow \\ \quad \quad \quad 4y^3 \\ \quad \quad \quad \downarrow \\ \quad \quad \quad g'(y) = 12y^2 \end{array} \right.$$

$$\nabla F(\omega) = \mathbb{E} [z z^T g'(\omega^T z)] - \lambda I$$

$$\approx \mathbb{E} [z z^T] \mathbb{E} [g'(\omega^T z)] - \lambda I$$

$$= \underbrace{\mathbb{E} [g'(\omega^T z)]}_{\text{SCALAR}} I - \lambda I$$

↑  
DIAGONAL MTX WITH IDENTICAL ELEMENTS

# Other Nonlinearities

**Newton method:**

$$\mathcal{W}(l_2+1) = \mathcal{W}(l_2) - \left[ \nabla F(\mathcal{W}(l_2)) \right]^{-1} F(\mathcal{W}(l_2))$$

$$= \mathcal{W}(l_2) - \frac{\mathbb{E}[z g(\mathcal{W}(l_2)^T z)] - \alpha \mathcal{W}(l_2)}{\mathbb{E}[g'(\mathcal{W}(l_2)^T z)] - \alpha}$$

$$\Rightarrow (\mathbb{E}[g'(\mathcal{W}(l_2)^T z)] - \alpha) \mathcal{W}(l_2+1) = (\mathbb{E}[g'(\mathcal{W}(l_2)^T z)] - \alpha) \mathcal{W}(l_2) + \alpha \mathcal{W}(l_2) - \mathbb{E}[z g(\mathcal{W}(l_2)^T z)]$$

$$\stackrel{C}{\uparrow}_{C \in \mathbb{R}} \mathcal{W}(l_2+1) = \mathbb{E}[z g(\mathcal{W}(l_2)^T z)] - \mathbb{E}[g'(\mathcal{W}(l_2)^T z)] \mathcal{W}(l_2)$$

**Algorithm:**

$$\Rightarrow \mathcal{W}(l_2+1) = \frac{\mathbb{E}[z g(\mathcal{W}(l_2)^T z)] - \mathbb{E}[g'(\mathcal{W}(l_2)^T z)] \mathcal{W}(l_2)}{\mathbb{E}[g'(\mathcal{W}(l_2)^T z)] - \alpha}$$

$$\mathcal{W}(l_2+1) = \mathcal{W}(l_2+1) / \|\mathcal{W}(l_2+1)\|$$

$$\left| \begin{array}{l} g_1(u) = \tanh(a_1 u) \\ g_2(u) = u \exp(-u^2/2) \end{array} \right.$$

# Fast ICA for several units

$$W = \begin{bmatrix} w_1^T \\ \vdots \\ w_n^T \end{bmatrix} \quad WW^T = I$$

WE NEED TO PREVENT DIFFERENT VECTORS CONVERGING TO THE SAME MAXIMA

IF WE ALREADY ESTIMATED  $w_1, \dots, w_p$   
THEN UPDATE  $w_{p+1}$

AND AFTER THAT

$$w_{p+1} = w_{p+1} - \sum_{j=1}^p (w_{p+1}^T w_j) w_j$$

GRAM-SCHMIDT  
DECORRELATION

OPTION 2:

$$W = (WW^T)^{-1/2} W$$

↑  
FROM SVD

$$WW^T = UDU^T \\ (WW^T)^{-1/2} = UD^{-1/2}U^T$$

SYMMETRIC  
DECORRELATION